

Hide-and-Seek: Forcing a Network to be Meticulous for Weakly-supervised Object and Action Localization

Krishna Kumar Singh and Yong Jae Lee
University of California, Davis

Abstract

We propose ‘Hide-and-Seek’, a weakly-supervised framework that aims to improve object localization in images and action localization in videos. Most existing weakly-supervised methods localize only the most discriminative parts of an object rather than all relevant parts, which leads to suboptimal performance. Our key idea is to hide patches in a training image randomly, forcing the network to seek other relevant parts when the most discriminative part is hidden. Our approach only needs to modify the input image and can work with any network designed for object localization. During testing, we do not need to hide any patches. Our Hide-and-Seek approach obtains superior performance compared to previous methods for weakly-supervised object localization on the ILSVRC dataset. We also demonstrate that our framework can be easily extended to weakly-supervised action localization.

1. Introduction

Weakly-supervised approaches have been proposed for various visual classification and localization tasks including object detection [54, 13, 9, 40, 3, 49, 42, 8, 31, 59, 39], semantic segmentation [32, 25] and visual attribute localization [2, 52, 55, 51, 38]. The main advantage of weakly-supervised learning is that it requires less detailed annotations compared to the fully-supervised setting, and therefore has the potential to use the vast weakly-annotated visual data available on the Web. For example, weakly-supervised object detectors can be trained using only image-level labels (‘dog’ or ‘no dog’) without any object location annotations.

Existing weakly-supervised methods identify discriminative patterns in the training data that frequently appear in one class and rarely in the remaining classes. This is done either explicitly by mining discriminative image regions or features [54, 13, 9, 40, 3, 41, 42, 8, 39] or implicitly by analyzing the higher-layer activation maps produced by a deep network trained for image classification [37, 31, 59]. However, due to intra-category variations or relying only on a classification objective, these methods often fail to identify the entire extent of the object and instead localize only the

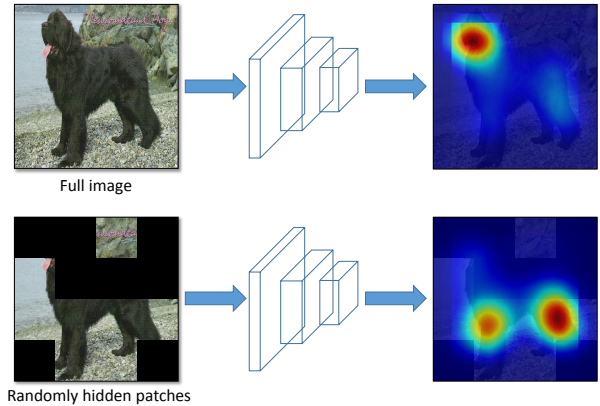


Figure 1. **Main idea.** (Top row) A network tends to focus on the most discriminative parts of an image (e.g., face of the dog) for classification. (Bottom row) By hiding images patches randomly, we can force the network to focus on other relevant object parts in order to correctly classify the image as ‘dog’.

most discriminative part.

Recent work tries to address this issue of identifying only the most discriminative part. Song et al. [42] combine multiple co-occurring discriminative regions to cover a larger extent of the object. While multiple selections ensure larger coverage, it does not guarantee selection of less discriminative patches of the object in the presence of many highly discriminative ones. Singh et al. [39] use motion cues and transfer tracked object boxes from weakly-labeled videos to the images. However, this approach requires additional weakly-labeled videos, which may not always be available. Finally, Zhou et al. [59] replace max pooling with global average pooling after the final convolution layer of an image classification network. Since average pooling aggregates activations across an entire feature map, it encourages the network to look beyond the most discriminative part (which would suffice for max pooling). However, the network can still avoid finding less discriminative parts if identifying a few highly-discriminative parts can lead to accurate classification performance, as shown in Figure 1 (top row).

Main Idea. In this paper, we take a very different approach to this problem. Instead of making algorithmic

changes [42, 59] or relying on external data [39], we make changes to the *input image*. The key idea is to *hide* patches from an image during training so that the model needs to *seek* the relevant object parts from what remains. We thus name our approach ‘Hide-and-Seek’. Figure 1 (bottom row) demonstrates the intuition: if we randomly remove some patches from the image then there is a possibility that the dog’s face, which is the most discriminative, will not be visible to the model. In this case, the model must seek other relevant parts like the tail and legs in order to do well on the classification task. By randomly hiding different patches in each training epoch, the model sees different parts of the image and is forced to focus on multiple relevant parts of the object beyond just the most discriminative one. Importantly, we only apply this random hiding of patches during training and not during testing. Since the full image is observed during testing, the data distribution will be different to that seen during training. We show that setting the hidden pixels’ value to be the data mean can allow the two distributions to match, and provide a theoretical justification.

Since Hide-and-Seek only alters the input image, it can easily be generalized to different neural networks and tasks. In this work, we demonstrate its applicability on AlexNet [27] and GoogLeNet [45], and apply the idea to weakly-supervised object localization in images and weakly-supervised action localization in videos. For the temporal action localization task (in which the start and end times of an action need to be found), random frame sequences are hidden while training a network on action classification, which forces the network to learn the relevant frames corresponding to an action.

Contributions. Our work has three main contributions: 1) We introduce the idea of Hide-and-Seek for weakly-supervised localization and produce state-of-the-art object localization results on the ILSVRC dataset [35]; 2) We demonstrate the generalizability of the approach on different networks and layers; 3) We extend the idea to the relatively unexplored task of weakly-supervised temporal action localization. We will share our Torch source code.

2. Related Work

Weakly-supervised object localization. Fully-supervised convolutional networks (CNNs) have demonstrated great performance on object detection [15, 14, 29], segmentation [30] and attribute localization [11, 58, 26], but require expensive human annotations for training (e.g., bounding box for object localization). To alleviate expensive annotation costs, weakly-supervised approaches learn using cheaper labels, for example, image-level labels for predicting an object’s location [54, 13, 9, 40, 3, 42, 49, 8, 31, 59].

Most weakly-supervised object localization approaches

mine discriminative features or patches in the data that frequently appear in one class and rarely in other classes [54, 13, 9, 40, 3, 7, 41, 42, 8]. However, these approaches tend to focus only on the most discriminative parts, and thus fail to cover the entire spatial extent of an object. In our approach, we hide image patches (randomly) during training, which forces our model to focus on multiple parts of an object and not just the most discriminative ones. Other methods use additional motion cues from weakly-labeled videos to improve object localization [34, 39]. While promising, such videos are not always readily available and can be challenging to obtain especially for static objects. In contrast, our method does not require any additional data or annotations.

Recent work modify CNN architectures designed for image classification so that the convolutional layers learn to localize objects while performing image classification [31, 59]. Other network architectures have been designed for weakly-supervised object detection [19, 4, 23]. Although these methods have significantly improved the state-of-the-art, they still essentially rely on a classification objective and thus can fail to capture the full extent of an object if the less discriminative parts do not help improve classification performance. We also rely on a classification objective. However, rather than modifying the CNN architecture, we instead modify the *input image* by hiding random patches from it. We demonstrate that this enforces the network to give attention to the less discriminative parts and ultimately localize a larger extent of the object.

Masking pixels or activations. Masking image patches has been applied for object localization [1], self-supervised feature learning [33], semantic segmentation [16, 10], generating hard occlusion training examples for object detection [53], and to visualize and understand what a CNN has learned [57]. In particular, for object localization, [57, 1] train a CNN for image classification and then localize the regions whose masking leads to a large drop in classification performance. Since these approaches mask out the image regions only during *testing* and not during training, the localized regions are limited to the highly-discriminative object parts. In our approach, image regions are masked during *training*, which enables the model to learn to focus on even the less discriminative object parts.

Dropout [43] and its variants [48, 46] are also related. There are two main differences: (1) these methods are designed to prevent overfitting while our work is designed to improve localization; and (2) in dropout, units in a layer are dropped randomly, while in our work, contiguous image regions or video frames are dropped. We demonstrate in the experiments that our approach produces significantly better localizations compared to dropout.

Action localization. Action localization is a well studied problem [28, 6, 50, 20, 22]. Recent CNN-based ap-

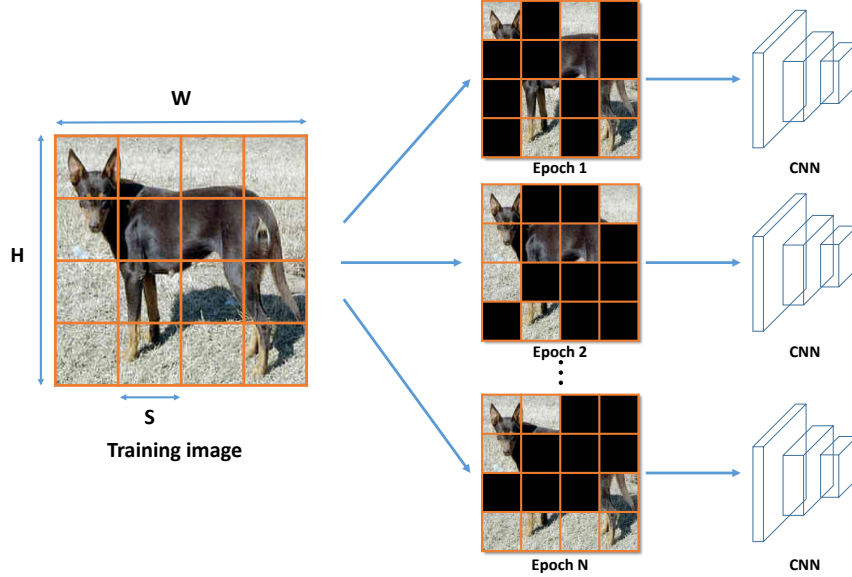


Figure 2. **Approach overview.** For each training image, we divide it into a grid of $S \times S$ patches. Each patch is then randomly hidden with probability p_{hide} and given as input to a CNN to learn image classification. The hidden patches change randomly across different epochs, which forces the network to focus on different parts of the object for learning image classification.

proaches [56, 36] have shown superior performance compared to previous hand-crafted approaches. These fully-supervised methods require the start and end time of an action in the video during the training to be annotated, which can be expensive to obtain. Weakly-supervised approaches learn from movie scripts [28, 12] or an ordered list of actions [5, 17]. Sun et al. [44] combine weakly-labeled videos with web images for action localization. In contrast to these approaches, our approach only uses a single video-level action label for temporal action localization.

3. Approach

In this section, we first describe our Hide-and-Seek algorithm for object localization in images followed by action localization in videos.

3.1. Weakly-supervised object localization

For weakly-supervised object localization, we are given a set of images $I_{set} = \{I_1, I_2, \dots, I_N\}$ in which each image I is labeled only with its category label. Our goal is to learn an object localizer that can predict both the category label as well as the bounding box for the object-of-interest in a new test image I_{test} . In order to learn the object localizer, we train a CNN which simultaneously learns to localize the object while performing the image classification task. While numerous approaches have been proposed to solve this problem, existing methods (e.g., [41, 8, 31, 59]) are prone to localizing only the most discriminative object parts, since those parts are sufficient for optimizing the classification task.

To enforce the network to learn all of the relevant parts

of an object, our key idea is to randomly hide patches of each input image I during training, as we explain next.

Hiding random image patches. The purpose of hiding patches is to show different parts of an object to the network while training it for the classification task. By hiding patches randomly, we can ensure that the most discriminative parts of an object are not always visible to the network, and thus *force* it to also focus on other relevant parts of the object. In this way, we can overcome the limitation of existing weakly-supervised methods that focus only on the most discriminative parts of an object.

Concretely, given a training image I of size $W \times H \times 3$, we first divide it into a grid with a fixed patch size of $S \times S \times 3$. This results in a total of $(W \times H)/(S \times S)$ patches. We then hide each patch with p_{hide} probability. For example, in Figure 2, the image is of size $224 \times 224 \times 3$, and it is divided into 16 patches of size $56 \times 56 \times 3$. Each patch is hidden with $p_{hide} = 0.5$ probability. We take the new image I' with the hidden patches, and feed it as a training input to a CNN for classification.

Importantly, for each image, we randomly hide a different set of patches. Also, for the same image, we randomly hide a different set of patches in each training epoch. This property allows the network to learn multiple relevant object parts for each image. For example, in Figure 2, the network sees a different I' in each epoch due to the randomness in hiding of the patches. In the first epoch, the dog’s face is hidden while its legs and tail are clearly visible. In contrast, in the second epoch, the face is visible while the legs and tail are hidden. Thus, the network is forced to learn all of the relevant parts of the dog rather than only the highly dis-

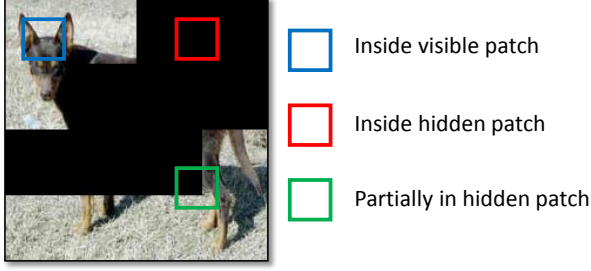


Figure 3. There are three types of convolutional filter activations after hiding patches: a convolution filter can be completely within a visible region (blue box), completely within a hidden region (red box), or partially within a visible/hidden region (green box).

crimative part (i.e., the face) in order to perform well in classifying the image as a ‘dog’.

We hide patches only during training. During testing, the full image, without any hidden patches, is given as input to the network. Since the network has learned to focus on multiple relevant parts during training, it is not necessary to hide any patches during testing. This is in direct contrast to [1], which hides patches during testing but not during training. For [1], since the network has already learned to focus on the most discriminative parts during training, it is essentially too late, and hiding patches during testing has no significant effect on localization performance.

Setting the hidden pixel values. There is an important detail that we must be careful about. Due to the discrepancy of hiding patches during training while not hiding patches during testing, the first convolutional layer activations during training versus testing will have different distributions. For a trained network to generalize well to new test data, the activation distributions should be roughly equal. That is, for any unit in a neural network that is connected to \mathbf{x} units with \mathbf{w} outgoing weights, the distribution of $\mathbf{w}^\top \mathbf{x}$ should be roughly the same during training and testing. However, in our setting, this will not necessarily be the case since some patches in each training image will be hidden while none of the patches in each test image will ever be hidden.

Specifically, in our setting, suppose that we have a convolution filter F with kernel size $K \times K$ and three-dimensional weights $W = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{k \times k}\}$, which is applied to an RGB patch $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{k \times k}\}$ in image I' . Denote \mathbf{v} as the vector representing the RGB value of every hidden pixel. There are three types of activations:

1. F is completely within a visible patch (Fig. 3, blue box). The corresponding output will be $\sum_{i=1}^{k \times k} \mathbf{w}_i^\top \mathbf{x}_i$.
2. F is completely within a hidden patch (Fig. 3, red box). The corresponding output will be $\sum_{i=1}^{k \times k} \mathbf{w}_i^\top \mathbf{v}$.
3. F is partially within a hidden patch (Fig. 3, green box). The corresponding output will be $\sum_{m \in \text{visible}} \mathbf{w}_m^\top \mathbf{x}_m + \sum_{n \in \text{hidden}} \mathbf{w}_n^\top \mathbf{v}$.

During testing, F will always be completely within a visible patch, and thus its output will be $\sum_{i=1}^{k \times k} \mathbf{w}_i^\top \mathbf{x}_i$. This matches the expected output during training in only the first case. For the remaining two cases, when F is completely or partially within a hidden patch, the activations will have a distribution that is different to those seen during the testing.

We resolve this issue by setting the RGB value \mathbf{v} of a hidden pixel to be equal to the mean RGB vector of the images over the entire dataset: $\mathbf{v} = \mu = \frac{1}{N_{\text{pixels}}} \sum_j \mathbf{x}_j$, where j indexes all pixels in the entire training dataset and N_{pixels} is the total number of pixels in the dataset. Why would this work? Essentially, we are assuming that in expectation, the output of a patch will be equal to that of an average-valued patch: $\mathbb{E}[\sum_{i=1}^{k \times k} \mathbf{w}_i^\top \mathbf{x}_i] = \sum_{i=1}^{k \times k} \mathbf{w}_i^\top \mu$. By replacing \mathbf{v} with μ , the outputs of both the second and third cases will be $\sum_{i=1}^{k \times k} \mathbf{w}_i^\top \mu$, and thus will match the expected output during testing (i.e., of a fully-visible patch).¹

This process is related to the scaling procedure in dropout [43], in which the outputs are scaled proportional to the drop rate during testing to match the expected output during training. In dropout, the outputs are dropped uniformly across the entire feature map, independently of spatial location. If we view our hiding of the patches as equivalent to “dropping” units, then in our case, we cannot have a global scale factor since the output of a patch depends on whether there are any hidden pixels. Thus, we instead set the hidden values to be the expected pixel value of the training data as described above, and do not scale the corresponding output. Empirically, we find that setting the hidden pixel in this way is important for the network to behave similarly during training and testing. If we set the hidden pixel to any other value, we observe a significant decrease in classification and localization performance.

Object localization network architecture. Our approach of hiding patches is independent of the network architecture and can be used with any CNN designed for object localization. For our experiments, we choose to use the network of Zhou et al. [59], which performs global average pooling (GAP) over the convolution feature maps to generate a class activation map (CAM) for the input image that represents the discriminative regions for a given class. This approach has shown state-of-the-art performance for the ILSVRC localization challenge [35] in the weakly-supervised setting, and existing CNN architectures like AlexNet [27] and GoogLeNet [45] can easily be modified to generate a CAM.

To generate a CAM for an image, global average pooling is performed after the last convolutional layer and the result is given to a classification layer to predict the image’s class probabilities. The weights associated with a class in

¹For the third case: $\sum_{m \in \text{visible}} \mathbf{w}_m^\top \mathbf{x}_m + \sum_{n \in \text{hidden}} \mathbf{w}_n^\top \mu \approx \sum_{m \in \text{visible}} \mathbf{w}_m^\top \mu + \sum_{n \in \text{hidden}} \mathbf{w}_n^\top \mu = \sum_{i=1}^{k \times k} \mathbf{w}_i^\top \mu$.

the classification layer represent the importance of the last convolutional layer’s feature maps for that class. More formally, denote $F = \{F_1, F_2, \dots, F_M\}$ to be the feature maps of the last convolutional layer and W as the $N \times M$ weight matrix of the classification layer, where N is number of classes. Then, the CAM for class c for image I is:

$$CAM(c, I) = \sum_{i=1}^M W(c, i) \cdot F_i(I). \quad (1)$$

Given the CAM for an image, we generate a bounding box using the method proposed in [59]. Briefly, we first threshold the CAM to produce a binary foreground/background map, and then find connected components among the foreground pixels. Finally, we fit a tight bounding box to the largest connected component. We refer the reader to [59] for more details.

Finally, our idea of hiding patches is not restricted to the input image layer. It can be also applied to the output of any convolutional layer in the same fashion, as we demonstrate in our experiments.

3.2. Weakly-supervised action localization

Given a set of untrimmed videos $V_{set} = \{V_1, V_2, \dots, V_N\}$ and video class labels, our goal here is to learn an action localizer that can predict the label of an action as well as its start and end time for a test video V_{test} . Again the key issue is that for any video, a network will focus mostly on the highly-discriminative frames in order to optimize classification accuracy instead of identifying all the relevant frames. Inspired by our idea of hiding patches in images, we propose to hide frames in videos to improve action localization.

Specifically, during training, we first uniformly sample F_{total} frames from each video. We then divide the F_{total} frames into continuous segments of fixed size $F_{segment}$; i.e., we have $F_{total}/F_{segment}$ segments. Just like with image patches, we hide each segment with probability p_{hide} before feeding it into a deep action localizer network. We generate class activation maps (CAM) using the procedure described in the previous section. In this case, our CAM is a one-dimensional map representing the discriminative frames for the action class. We apply thresholding on this map to obtain the start and end times for the action class.

4. Experiments

In this section, we perform quantitative and qualitative evaluations of our Hide-and-Seek method for object localization in images and action localization in videos. In addition, we perform ablative studies to compare the different design choices of our algorithm.

Datasets and evaluation metrics. We use ILSVRC 2016 [35] to evaluate object localization accuracy. For

training, we use 1.2 million images with their class labels (1000 categories). We compare our approach with the baselines on both the validation and test data. We use two evaluation metrics to measure performance: 1) Top-1 localization accuracy (*Top-1 Loc*): fraction of images for which the predicted class with the highest probability is the same as the ground-truth class *and* the predicted bounding box for that class has more than 50% IoU with the ground-truth box; 2) Localization accuracy with known ground-truth class (*GT-known Loc*): fraction of images for which the predicted bounding box for the ground-truth class has more than 50% IoU with the ground-truth box. As our approach is primarily designed to improve localization accuracy, we use the second evaluation criterion to measure localization accuracy independent from classification performance.

For action localization, we use the THUMOS 2014 validation data [21], which consists of 1010 untrimmed videos belonging to 101 action classes. We train our network over all untrimmed videos for the classification task and then evaluate localization on the 20 classes that have temporal annotations. Each video can contain multiple instances of a class. For evaluation we compute mean average precision (mAP), and consider a prediction to be correct if it has IoU $> \theta$ with the ground-truth. We vary θ to be 0.1, 0.2, 0.3, 0.4, and 0.5. As we are focusing on localization ability of the network, we assume that we know the ground-truth class label of the video.

Implementation details. To learn the object localizer, we use the same modified AlexNet and GoogLeNet networks introduced in Zhou et al. [59] (AlexNet-GAP and GoogLeNet-GAP). AlexNet-GAP is identical to AlexNet until pool5 after which two new convolutional layers are added. Similarly for GoogLeNet-GAP, layers after inception-4e are removed and a single convolutional layer is added. For both AlexNet-GAP and GoogLeNet-GAP, the output of the last convolutional layer goes to a global average pooling (GAP) layer, followed by a softmax layer for classification. Each added convolutional layer has 512 and 1024 kernels of size 3×3 , stride 1, and pad 1 for AlexNet-GAP and GoogLeNet-GAP, respectively.

We implement the networks using Torch. We train the networks from scratch for 55 and 40 epochs for AlexNet-GAP and GoogLeNet-GAP, respectively, with a batch size of 128 and an initial learning rate of 0.01. We decrease the learning rate by $0.1 \times$ after every 25 epochs. We add a batch normalization [18] layer after every convolutional layer to help convergence of GoogLeNet-GAP. For simplicity, unlike the original AlexNet architecture [27], we do not group the convolutional filters together (it produces statistically the same *Top-1 Loc* accuracy as the grouped version for both AlexNet-GAP and GoogLeNet-GAP). The network remains exactly the same with (during training) and without (during testing) hidden patches. To obtain the binary fore-

Methods	GT-known Loc	Top-1 Loc
AlexNet-GAP [59]	54.99 ²	36.25
AlexNet-HaS-16	58.00	36.89
AlexNet-HaS-32	58.92	37.46
AlexNet-HaS-44	58.74	37.71
AlexNet-HaS-56	58.62	37.50
AlexNet-HaS-Mixed	59.36	37.64
GoogLeNet-GAP [59]	58.66 ²	43.60
GoogLeNet-HaS-16	60.10	44.87
GoogLeNet-HaS-32	60.57	45.47
GoogLeNet-HaS-44	60.39	45.01
GoogLeNet-HaS-56	60.22	45.03

Table 1. Localization accuracy on ILSVRC validation data with different patch sizes for hiding. Our Hide-and-Seek always performs better than AlexNet-GAP/GoogLeNet-GAP [59], which sees the full image.

ground/background map, 20% and 30% of the max value of the CAM is chosen as the threshold for the AlexNet-GAP and GoogLeNet-GAP, respectively. These thresholds were chosen by observing a few qualitative results on the training data. During testing, we average 10 crops (4 corners plus center, and same with horizontal flip) to obtain class probabilities and localization maps. We find similar localization/classification performance when fine-tuning pre-trained networks.

For action localization, we compute C3D [47] fc7 features using a model pre-trained on Sports 1 million [24]. We compute 10 features/second (each feature is computed over 16 frames) and uniformly sample 2000 features from the video. We then divide the video into 20 equal-length video segments where each segment consists of $F_{segment} = 100$ features. During training, we hide each segment with $p_{hide} = 0.5$ probability. For action classification, we feed the C3D features as input to a CNN with two convolution layers followed by a global max pooling and softmax classification layer. Each conv layer has 500 kernels of size 1×1 , stride 1. For any hidden frame, we assign it the mean C3D feature computed over the dataset. For thresholding, 50% of the max value of the CAM is chosen. All continuous segments after thresholding are considered predictions.

4.1. Object localization quantitative results

We first analyze object localization accuracy on the ILSVRC validation data. Table 1 shows the results using the *Top-1 Loc* and *GT-known Loc* evaluation metrics. AlexNet-GAP [59] is our baseline in which the network has seen the full image during training without any hidden patches. Alex-HaS-N is our approach, in which patches of size $N \times N$ are hidden with 0.5 probability during training.

Which patch size N should we choose? We explored four different patch sizes $N = \{16, 32, 44, 56\}$, and each performs significantly better than AlexNet-GAP for both

²[59] does not provide GT-known loc, so we compute on our own GAP implementations, which achieve similar Top-1 Loc accuracy.

Methods	GT-known Loc	Top-1 Loc
Backprop on AlexNet [37]	-	34.83
AlexNet-GAP [59]	54.99	36.25
Ours	58.74	37.71
AlexNet-GAP-ensemble	57.02	38.69
Ours-ensemble	60.33	40.57
Backprop on GoogLeNet [37]	-	38.69
GoogLeNet-GAP [59]	58.66	43.60
Ours	60.57	45.47

Table 2. Localization accuracy on ILSVRC val data compared to state-of-the-art. Our method outperforms all previous methods.

GT-known Loc and *Top-1 Loc*. Our GoogLeNet-HaS-N models also outperforms GoogLeNet-GAP for all patch sizes. These results clearly show that hiding patches during training leads to better localization. Although our approach might lose some classification accuracy since it has never seen a complete image and thus may not have learned to relate certain parts, the huge boost in localization performance (which can be seen by comparing the *GT-known Loc* accuracies) makes up for any potential classification loss.

We also train a network (AlexNet-HaS-Mixed) with mixed patch sizes. During training, for each image in every epoch, the patch size N to hide is chosen randomly from 16, 32, 44 and 56. Since different sized patches are hidden, the network can learn complementary information about different parts of an object (e.g., small/large patches are more suitable to hide smaller/larger parts). Indeed, our results show that we do better for *GT-known Loc* using AlexNet-HaS-Mixed.

Comparison to state-of-the-art. Next, we choose our best model for AlexNet and GoogLeNet, and compare it with state-of-the-art methods on ILSVRC validation data; see Table 2. Our method performs 3.75% and 1.46% points better than AlexNet-GAP [59] on *GT-known Loc* and *Top-1 Loc*, respectively. For GoogLeNet, our model gets a boost of 1.91% and 1.87% points compared to GoogLeNet-GAP for *GT-known Loc* and *Top-1 Loc* accuracy, respectively. Importantly, these gains are obtained simply by changing the input image without changing the network architecture.

Ensemble model. Since each patch size provides complementary information (as seen in the previous section), we also create an ensemble model of different patch sizes (Ours-ensemble). To produce the final localization for an image, we average the CAMs obtained using AlexNet-HaS-16, 32, 44, and 56, while for classification, we average the classification probabilities of all four models as well as the probability obtained using AlexNet-GAP. Although AlexNet-GAP has been trained using the full image, which is not optimal for localization, it is still well-suited for classification. This ensemble model gives a boost of 5.34% and 4.32% over AlexNet-GAP for *GT-known Loc* and *Top-1 Loc*, respectively. For a more fair comparison, we also combine the results of five independent AlexNet-GAPs to

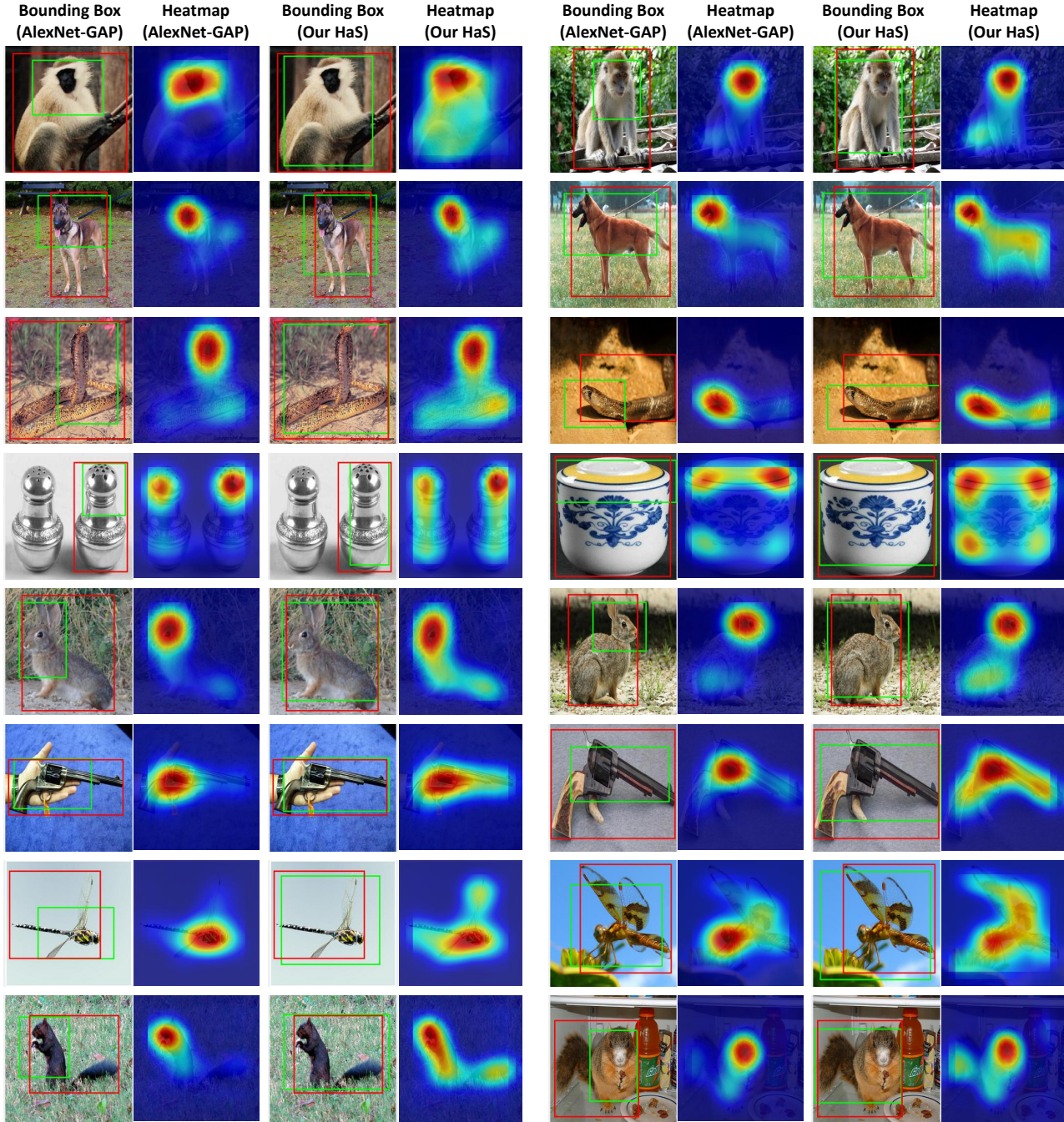


Figure 4. Qualitative object localization results. We compare our approach with AlexNet-GAP [59] on the ILVRC validation data. For each image, we show the bounding box and CAM obtained by AlexNet-GAP (left) and our method (right). Our Hide-and-Seek approach localizes multiple relevant parts of an object whereas AlexNet-GAP mainly focuses only on the most discriminative part. For example, in the first, second and fifth rows, our method localizes the full body of the animals while AlexNet-GAP only focuses on the face. Similarly, our method can capture the tail of the squirrels and snakes in the third and last rows, which are missed by AlexNet-GAP. We can even localize the wings of an insect, and front part of the gun in the second and third last rows, respectively.

create an ensemble baseline. Ours-ensemble outperforms this strong baseline (AlexNet-GAP-ensemble) by 3.31% and 1.88% for *GT-known Loc* and *Top-1 Loc*, respectively.

4.2. Object localization qualitative results

In Fig. 4, we visualize the class activation map (CAM) and bounding box obtained by our approach (implementation of HaS using the AlexNet-GAP network) versus those

Methods	GT-known Loc	Top-1 Loc
Ours	58.74	37.71
AlexNet-dropout-trainonly	42.20	07.68
AlexNet-dropout-traintest	53.55	31.72

Table 3. Comparison with Dropout [43]. Our approach performs better for localization.

Methods	GT-known Loc	Top-1 Loc
AlexNet-GAP	54.99	35.70
AlexNet-Avg-HaS	58.62	37.50
AlexNet-GMP	50.42	32.55
AlexNet-Max-HaS	59.48	37.74

Table 4. Global average pooling (GAP) vs. global max pooling (GMP). Unlike [59], for Hide-and-Seek GMP still performs well for localization. For this experiment, we use patch size 56.

obtained with AlexNet-GAP. In each image pair, the first image shows the predicted (green) and ground-truth (red) bounding box. The second image shows the CAM, i.e., where the network is focusing for that class.

Our approach localizes more relevant parts of an object compared to AlexNet-GAP and is not confined to only the most discriminative part. For example, in the first, second, and fifth rows AlexNet-GAP only focuses on the face of the animals, whereas our approach also localizes parts of the body. Similarly, in the third and last rows AlexNet-GAP misses the tail for the snake and squirrel while our approach gets the tail.

4.3. Further Analysis of Hide-and-Seek

Comparison with dropout. Dropout [43] has been extensively used to reduce overfitting in deep networks. Although it is not designed to improve localization, the dropping of units is related to our hiding of patches. We therefore conduct an experiment in which 50% dropout is applied at the image layer. We noticed that due to the large dropout rate at the pixel-level, the learned filters develop a bias towards a dropped-out version of the images and produces significantly inferior classification and localization performance, even with the proper scaling (AlexNet-dropout-trainonly). If we also do dropout during testing (AlexNet-dropout-traintest) then performance improves but is still much lower compared to our approach (Table 3). Since dropout drops pixels (and RGB channels) randomly, information from the most relevant parts of an object will still be seen by the network with high probability, which makes it likely to focus on only the most discriminative parts.

Do we need global average pooling? [59] showed that GAP is better than global max pooling (GMP) for object localization, since average pooling encourages the network to focus on all the discriminative parts. For max pooling, only the most discriminative parts need to contribute. But is global max pooling hopeless for localization?

With our Hide-and-Seek, even with max pooling, the

Methods	GT-known Loc	Top-1 Loc
AlexNet-GAP	54.99	35.70
AlexNet-HaS-conv1-5	57.49	37.02
AlexNet-HaS-conv1-11	58.49	37.52

Table 5. Applying Hide-and-Seek to the first convolutional layer. The improvement in localization accuracy over [59] shows the generality of the idea.

Methods	GT-known Loc	Top-1 Loc
AlexNet-HaS-25%	57.64	37.89
AlexNet-HaS-33%	58.27	38.19
AlexNet-HaS-50%	58.62	37.38
AlexNet-HaS-66%	58.76	35.90
AlexNet-HaS-75%	58.52	34.38

Table 6. Varying the probability of hiding patches. Higher probabilities lead to decrease in *Top-1 Loc* whereas lower probability leads to smaller *GT-known Loc*. For this experiment, we use patch size 56.

network is forced to focus on a different discriminative part in each training epoch, since different parts of the image will be hidden (i.e., it cannot keep selecting the single most discriminative part). In Table 4, we see that max pooling (AlexNet-GMP) is inferior to average pooling (AlexNet-GAP) for the baselines. But with Hide-and-Seek, max pooling (AlexNet-Max-HaS) localization accuracy increases by a big margin and even slightly outperforms average pooling (AlexNet-Avg-HaS). The slight improvement over average pooling is likely due to max pooling being more robust to noise.

Hide-and-Seek in convolutional layers. We next apply our idea to convolutional layers. We divide the convolutional feature maps into a grid and hide each patch (and all of its corresponding channels) with 0.5 probability. We hide patches of size 5 (AlexNet-HaS-conv1-5) and 11 (AlexNet-HaS-conv1-11) in the conv1 feature map (which has size $55 \times 55 \times 96$). Table 5 shows that this leads to a big boost in performance compared to the baseline AlexNet-GAP. This shows that our idea of randomly hiding patches can be generalized to the convolutional layers.

Probability of hiding. For all previous experiments, we hid the patches with 50% probability. In Table 6, we measure the *GT-known Loc* and *Top-1 Loc* when we use different hiding probabilities. If we increase the probability then *GT-known Loc* remains almost the same while *Top-1 Loc* decreases a lot. This happens because the network sees fewer pixels when the hiding probability is high; as a result, classification accuracy reduces and *Top-1 Loc* drops. If we decrease the probability then *GT-known Loc* decreases but our *Top-1 Loc* improves. In this case, the network sees more pixels so its classification improves but since less parts are hidden, it will focus more on only the discriminative parts decreasing its localization ability.

Methods	IOU thresh = 0.1	0.2	0.3	0.4	0.5
Video-full	34.23	25.68	17.72	11.00	6.11
Video-HaS	36.44	27.84	19.49	12.66	6.84

Table 7. Action localization accuracy on THUMOS validation data. Across all 5 IoU thresholds, our Video-HaS outperforms the full video baseline (Video-full).

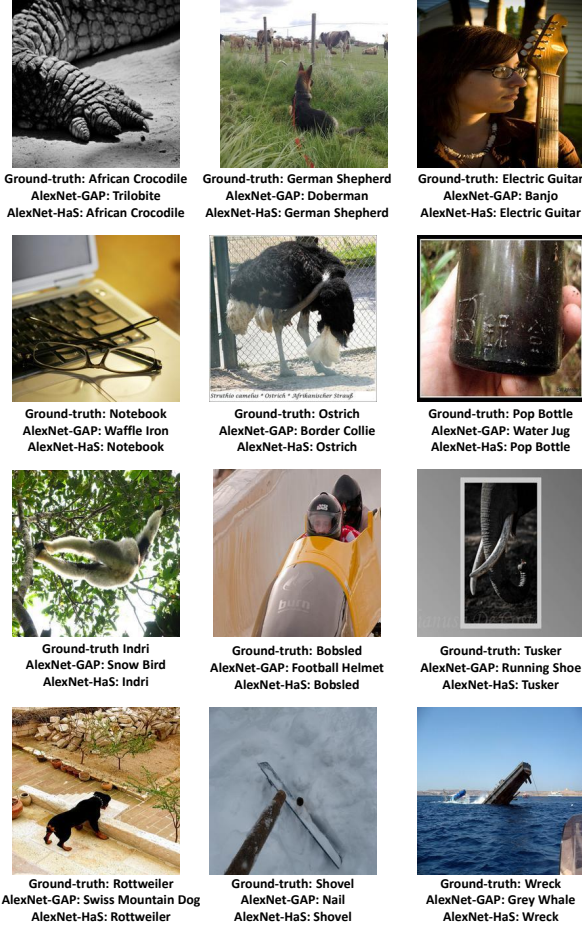


Figure 5. Comparison of our AlexNet-HaS vs. the AlexNet-GAP baseline for classification of challenging images. For each image, we show the ground-truth label followed by top class predicted by AlexNet-GAP and AlexNet-HaS. AlexNet-HaS is able to classify the images correctly even when they are partially occluded (African crocodile, electric guitar, notebook, pop bottle, bobsled, tusker, shovel and wreck). Even when the most discriminative part is hidden, our AlexNet-HaS classifies the image correctly; for example, the faces of the German Shepherd, ostrich, indri and rottweiler are hidden but our AlexNet-HaS is still able to classify them correctly.

4.4. Classification of challenging images

In our Hide-and-Seek (HaS) approach, the network is trained using images in which patches are hidden randomly. This gives the network the ability to classify images correctly even when the objects are partially-occluded and when its most discriminative parts are not visible. In Fig-

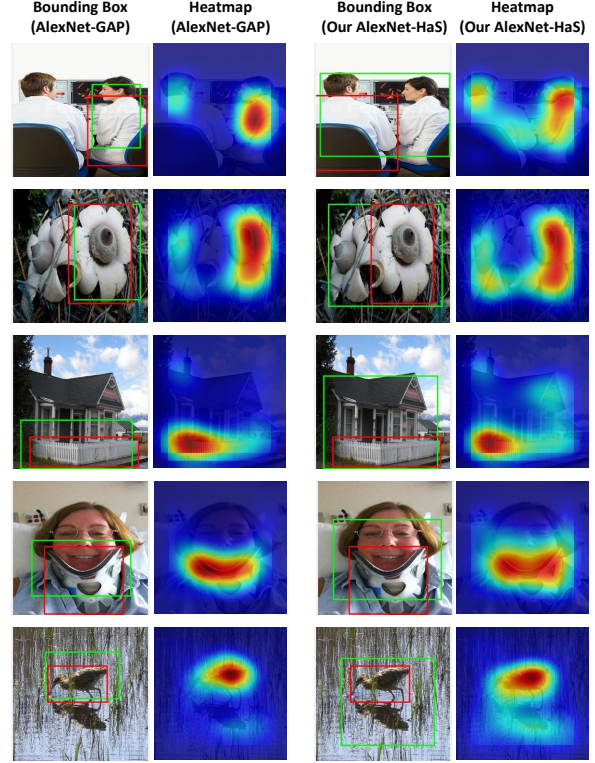


Figure 6. Example failure cases of AlexNet-HaS. For each image we show the bounding box (red: ground-truth, green: predicted) and CAM obtained by the AlexNet-GAP baseline (left) and our AlexNet-HaS approach (right). In the first two rows, our method fails due to merging of the localization of multiple instances of the object-of-interest. In the third and fourth rows, it fails due to strong co-occurrence of contextual objects with the object-of-interest. In the last row our localizer gets confused due to the reflection of the bird.

ure 5, we show challenging cases for which AlexNet-GAP fails but our AlexNet-HaS successfully classifies the images. Our AlexNet-HaS can correctly classify ‘African Crocodile’ and ‘Notebook’ by just looking at the leg and keypad, respectively. It can also classify ‘German Shepherd’, ‘Ostrich’, ‘Indri’ and ‘Rottweiler’ correctly without looking at the face, which is the most discriminative part.

4.5. Failure cases of Hide-and-Seek

Both the quantitative (Table 2) and qualitative results (Figure 4) in the paper show that overall our Hide-and-see (HaS) approach leads to better localization compared to the GAP baseline. Still, HaS is not perfect and there are some typical scenarios where it fails and produces inferior localization compared to GAP.

Figure 6 shows example failure cases of AlexNet-HaS compared to AlexNet-GAP. In the first two rows, HaS fails to localize a single object instance because there are multiple instances of the same object that are spatially close to

each other. This leads to our approach merging the localizations of the two object instances together. For example in the first row, our localization of the two lab coats are merged together to produce a bigger bounding box containing both of them. In contrast, AlexNet-GAP produces a more selective localization (focusing mainly on only the lab coat on the right), which leads to a bounding box that covers only a single lab coat. In the third and fourth rows, failure occurs due to the strong co-occurrence of the contextual objects near the object-of-interest. Specifically, in the third row, our AlexNet-HaS localizes parts of the house (context) along with the fence (object-of-interest) because house co-occurs with fences frequently. As a result, when parts of the fence are hidden during training the network starts to focus on the house regions in order to do well for the fence classification task. Finally, in the last row, our AlexNet-HaS localizes both the bird and its reflection in the water, which leads to an incorrect bounding box.

4.6. Action localization results

Finally, we evaluate action localization accuracy. We compare our approach (Video-HaS), which randomly hides frame segments while learning action classification, with a baseline that sees the full video (Video-full). Table 7 shows the result on THUMOS validation data. Video-HaS consistently outperforms Video-full, which shows that hiding frames forces our network to focus on more relevant frames, which ultimately leads to better action localization.

In Figure 7, we compare the temporal action localization results of our approach of randomly hiding frame segments while learning an action classifier (Video-HaS) versus the baseline approach of showing the whole video during training (Video-full). For each action, we uniformly sample the frames and show: 1) Ground-truth (first row, frames belonging to action have green boundary), 2) Video-full (second row, localized frames have red boundary) and 3) Video-HaS (third row, localized frames have red boundary).

From Figure 7, we can see that our Video-HaS localizes most of the temporal extent of an action while Video-full only localizes some key moments. For example, in the case of javelin throw (second example), Video-HaS localizes all the frames associated with the action where as Video-full only localizes a frame in which the javelin is thrown. In the third example, Video-full localizes only the beginning part of high jump while Video-HaS localizes all relevant frames. In last row, we show a failure case of Video-HaS in which it incorrectly localizes beyond the temporal extent of diving. Since frames containing a swimming pool follow the diving action frequently, when the diving frames are hidden the network starts focusing on the context frames containing swimming pool to classify the action as diving.

5. Conclusion

We presented ‘Hide-and-Seek’, a novel weakly-supervised framework to improve object localization in images and temporal action localization in videos. By randomly hiding patches/frames in a training image/video, we force the network to learn to focus on multiple relevant parts of an object/action. Our extensive experiments showed improved localization accuracy over state-of-the-art methods.

References

- [1] L. Bazzani, B. A., D. Anguelov, and L. Torresani. Self-taught object localization with deep networks. In *WACV*, 2016. 2, 4
- [2] T. Berg, A. Berg, and J. Shih. Automatic attribute discovery and characterization from noisy web data. In *ECCV*, 2010. 1
- [3] H. Bilen, M. Pedersoli, and T. Tuytelaars. Weakly supervised object detection with posterior regularization. In *BMVC*, 2014. 1, 2
- [4] H. Bilen and A. Vedaldi. Weakly supervised deep detection networks. In *CVPR*, 2016. 2
- [5] P. Bojanowski, R. Lajugie, F. Bach, I. Laptev, J. Ponce, C. Schmid, and J. Sivic. Weakly supervised action labeling in videos under ordering constraints. In *ECCV*, 2014. 3
- [6] C. Y. Chen and K. Grauman. Efficient activity detection with max-subgraph search. In *CVPR*, 2012. 2
- [7] R. Cinbis, J. Verbeek, and C. Schmid. Multi-fold MIL Training for Weakly Supervised Object Localization. In *CVPR*, 2014. 2
- [8] R. Cinbis, J. Verbeek, and C. Schmid. Weakly supervised object localization with multi-fold multiple instance learning. In *arXiv:1503.00949*, 2015. 1, 2, 3
- [9] D. J. Crandall and D. P. Huttenlocher. Weakly supervised learning of part-based spatial models for visual object recognition. In *ECCV*, 2006. 1, 2
- [10] J. Dai, K. He, and J. Sun. Convolutional feature masking for joint object and stuff segmentation. In *CVPR*, 2015. 2
- [11] K. Duan, D. Parikh, D. Crandall, and K. Grauman. Discovering localized attributes for fine-grained recognition. In *CVPR*, 2012. 2
- [12] O. Duchenne, I. Laptev, J. Sivic, F. Bach, and J. Ponce. Automatic annotation of human actions in video. In *ICCV*, 2009. 3
- [13] R. Fergus, P. Perona, and A. Zisserman. Object Class Recognition by Unsupervised Scale-Invariant Learning. In *CVPR*, 2003. 1, 2
- [14] R. Girshick. Fast r-cnn. In *ICCV*, 2015. 2
- [15] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. In *CVPR*, 2014. 2
- [16] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik. Simultaneous detection and segmentation. In *ECCV*, 2014. 2
- [17] D.-A. Huang, L. Fei-Fei, and J. C. Nibbles. Connectionist temporal modeling for weakly supervised action labeling. In *ECCV*, 2016. 3
- [18] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 5

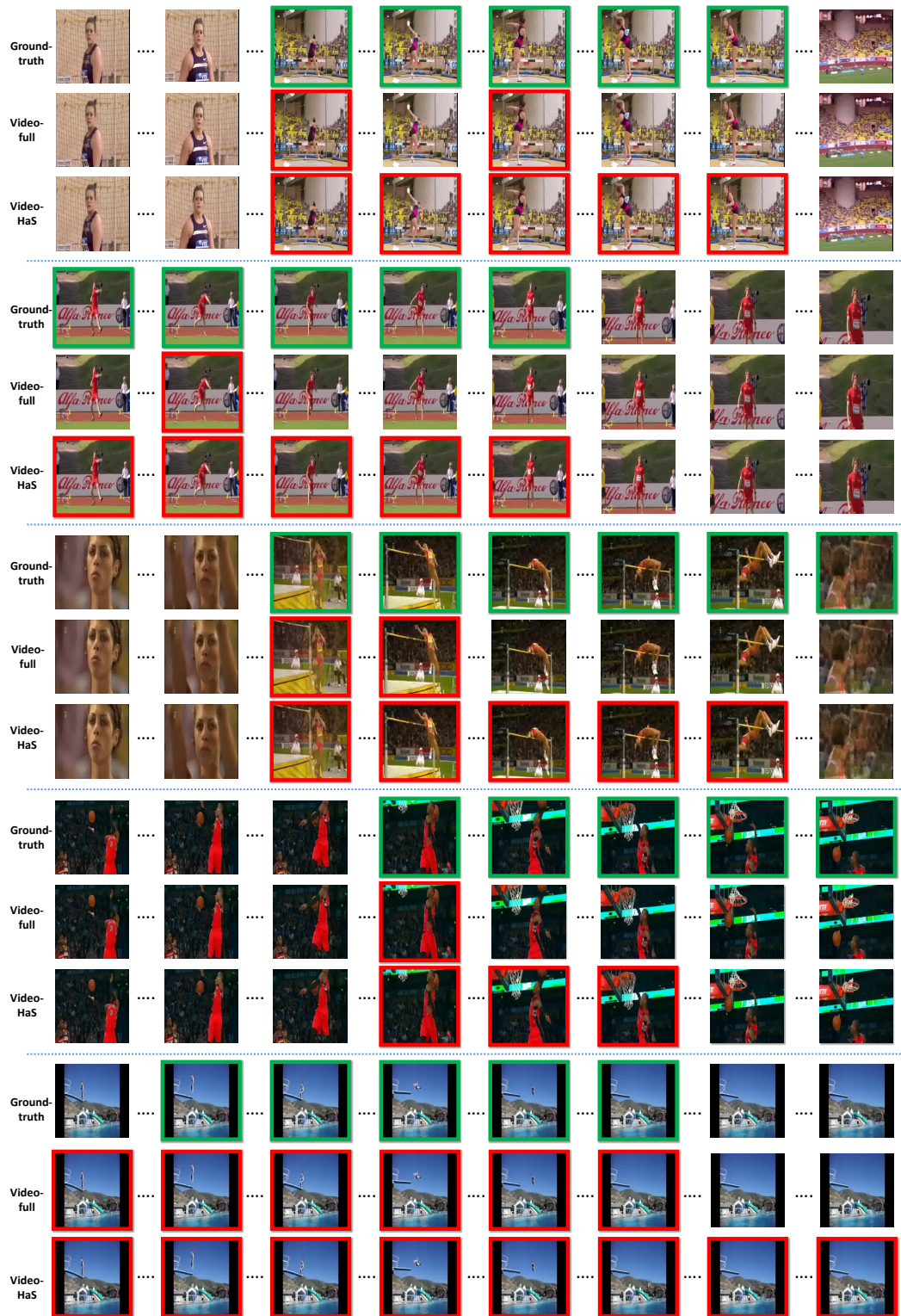


Figure 7. Comparison of action localization between the Video-full baseline and our method of Video-HaS. For each action, we uniformly sample the frames and show the ground-truth in the first row (frames with a green boundary belong to the action), followed by the Video-full and Video-HaS localizations (frames with a red boundary). For each action (except the last one), Video-HaS localizes the full extent of the action more accurately compared to Video-full, which tends to localize only some key frames. For example in the third example, Video-full only localizes the initial part of high-jump whereas Video-HaS localizes all relevant frames. In the last example, we show a failure case of our Video-HaS, in which it incorrectly localizes the last two frames as diving due to the co-occurring swimming pool context.

- [19] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu. Spatial transformer networks. In *NIPS*, 2015. 2
- [20] H. Jhuang, J. Gall, S. Zuffi, C. Schmid, and M. J. Black. Towards understanding action recognition. In *ICCV*, 2013. 2
- [21] Y.-G. Jiang, J. Liu, A. Roshan Zamir, G. Toderici, I. Laptev, M. Shah, and R. Sukthankar. THUMOS challenge: Action recognition with a large number of classes. <http://csrcv.ucf.edu/THUMOS14/>, 2014. 5
- [22] V. Kantorov and I. Laptev. Efficient feature extraction, encoding and classification for action recognition. In *CVPR*, 2014. 2
- [23] V. Kantorov, M. Oquab, M. Cho, and I. Laptev. Contextlocnet: Context-aware deep network models for weakly supervised localization. In *ECCV*, 2016. 2
- [24] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, 2014. 6
- [25] A. Khoreva, R. Benenson, M. Omran, M. Hein, and B. Schiele. Weakly supervised object boundaries. In *CVPR*, 2016. 1
- [26] M. Kiapour, K. Yamaguchi, A. C. Berg, and T. L. Berg. Hipster wars: Discovering elements of fashion styles. In *ECCV*, 2014. 2
- [27] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet Classification with Deep Convolutional Neural Networks. In *NIPS*, 2012. 2, 4, 5
- [28] I. Laptev, M. Marszalek, C. Schmid, and B. Rozenfeld. Learning realistic human actions from movies. In *CVPR*, 2008. 2, 3
- [29] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *ECCV*, 2016. 2
- [30] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 2
- [31] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Is object localization for free? weakly-supervised learning with convolutional neural networks. In *CVPR*, 2015. 1, 2, 3
- [32] D. Pathak, P. Krähenbühl, and T. Darrell. Constrained convolutional neural networks for weakly supervised segmentation. In *ICCV*, 2015. 1
- [33] D. Pathak, P. Krähenbühl, J. Donahue, T. Darrell, and A. Efros. Context encoders: Feature learning by inpainting. In *CVPR*, 2016. 2
- [34] A. Prest, C. Leistner, J. Civera, C. Schmid, and V. Ferrari. Learning Object Class Detectors from Weakly Annotated Video. In *CVPR*, 2012. 2
- [35] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 2015. 2, 4, 5
- [36] Z. Shou, D. Wang, and S.-F. Chang. Temporal action localization in untrimmed videos via multi-stage cnns. In *CVPR*, 2016. 2
- [37] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *ICLR Workshop*, 2014. 1, 6
- [38] K. K. Singh and Y. J. Lee. End-to-end localization and ranking for relative attributes. In *ECCV*, 2016. 1
- [39] K. K. Singh, F. Xiao, and Y. J. Lee. Track and transfer: Watching videos to simulate strong human supervision for weakly-supervised object detection. In *CVPR*, 2016. 1, 2
- [40] P. Siva, C. Russell, and T. Xiang. In Defence of Negative Mining for Annotating Weakly Labelled Data. In *ECCV*, 2012. 1, 2
- [41] H. O. Song, R. Girshick, S. Jegelka, J. Mairal, Z. Harchaoui, and T. Darrell. On Learning to Localize Objects with Minimal Supervision. In *ICML*, 2014. 1, 2, 3
- [42] H. O. Song, Y. J. Lee, S. Jegelka, and T. Darrell. Weakly-supervised discovery of visual pattern configurations. In *NIPS*, 2014. 1, 2
- [43] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *JMLR*, 2014. 2, 4, 8
- [44] C. Sun, S. Shetty, R. Sukthankar, and R. Nevatia. Temporal localization of fine-grained actions in videos by domain transfer from web images. In *ACM Multimedia*, 2015. 3
- [45] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 2, 4
- [46] J. Tompson, R. Goroshin, A. Jain, Y. LeCun, and C. Bregler. Efficient object localization using convolutional networks. In *CVPR*, 2015. 2
- [47] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*, 2015. 6
- [48] L. Wan, M. Zeiler, S. Zhang, Y. LeCun, and R. Fergus. Regularization of neural network using dropconnect. In *ICML*, 2013. 2
- [49] C. Wang, W. Ren, K. Huang, and T. Tan. Weakly supervised object localization with latent category learning. In *ECCV*, 2014. 1, 2
- [50] H. Wang and C. Schmid. Action recognition with improved trajectories. In *ICCV*, 2013. 2
- [51] J. Wang, Y. Cheng, and R. Schmidt Feris. Walk and learn: Facial attribute representation learning from egocentric video and contextual data. In *CVPR*, 2016. 1
- [52] S. Wang, J. Joo, Y. Wang, and S. C. Zhu. Weakly supervised learning for attribute localization in outdoor scenes. In *CVPR*, 2013. 1
- [53] X. Wang, A. Shrivastava, and A. Gupta. A-fast-rcnn: Hard positive generation via adversary for object detection. In *arXiv:1704.03414*, 2017. 2
- [54] M. Weber, M. Welling, and P. Perona. Unsupervised Learning of Models for Recognition. In *ECCV*, 2000. 1, 2
- [55] F. Xiao and Y. J. Lee. Discovering the spatial extent of relative attributes. In *ICCV*, 2015. 1
- [56] S. Yeung, O. Russakovsky, G. Mori, and L. Fei-Fei. End-to-end learning of action detection from frame glimpses in videos. In *CVPR*, 2016. 2
- [57] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *ECCV*, 2014. 2
- [58] N. Zhang, M. Paluri, M. Ranzato, T. Darrell, and L. Bourdev. PANDA: Pose Aligned Networks for Deep Attribute Modeling. In *CVPR*, 2014. 2

- [59] B. Zhou, A. Khosla, L. A., A. Oliva, and A. Torralba. Learning deep features for discriminative localization. In *CVPR*, 2016. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#)